

20.4 Huffman Coding and Compression of Data

A lossless data compression algorithm takes a string of symbols (typically ASCII characters or bytes) and translates it *reversibly* into another string, one that is *on the average* of shorter length. The words “on the average” are crucial; it is obvious that no reversible algorithm can make all strings shorter — there just aren’t enough short strings to be in one-to-one correspondence with longer strings. Compression algorithms are possible only when, on the input side, some strings, or some input symbols, are more common than others. These can then be encoded in fewer bits than rarer input strings or symbols, giving a net average gain.

There exist many, quite different, compression techniques, corresponding to different ways of detecting and using departures from equiprobability in input strings. In this section and the next we shall consider only *variable length codes* with *defined word* inputs. In these, the input is sliced into fixed units, for example ASCII characters, while the corresponding output comes in chunks of variable size. The simplest such method is Huffman coding [1], discussed in this section. Another example, *arithmetic compression*, is discussed in §20.5.

At the opposite extreme from defined-word, variable length codes are schemes that divide up the *input* into units of variable length (words or phrases of English text, for example) and then transmit these, often with a fixed-length output code. The most widely used code of this type is the Ziv-Lempel code [2]. References [3-6] give the flavor of some other compression techniques, with references to the large literature.

The idea behind Huffman coding is simply to use shorter bit patterns for more common characters. We can make this idea quantitative by considering the concept of *entropy*. Suppose the input alphabet has N_{ch} characters, and that these occur in the input string with respective probabilities p_i , $i = 1, \dots, N_{ch}$, so that $\sum p_i = 1$. Then the fundamental theorem of information theory says that strings consisting of independently random sequences of these characters (a conservative, but not always realistic assumption) require, on the average, at least

$$H = - \sum p_i \log_2 p_i \quad (20.4.1)$$

bits per character. Here H is the entropy of the probability distribution. Moreover, coding schemes exist which approach the bound arbitrarily closely. For the case of equiprobable characters, with all $p_i = 1/N_{ch}$, one easily sees that $H = \log_2 N_{ch}$, which is the case of no compression at all. Any other set of p_i ’s gives a smaller entropy, allowing some useful compression.

Notice that the bound of (20.4.1) would be achieved if we could encode character i with a code of length $L_i = -\log_2 p_i$ bits: Equation (20.4.1) would then be the average $\sum p_i L_i$. The trouble with such a scheme is that $-\log_2 p_i$ is not generally an integer. How can we encode the letter “Q” in 5.32 bits? Huffman coding makes a stab at this by, in effect, approximating all the probabilities p_i by integer powers of $1/2$, so that all the L_i ’s are integral. If all the p_i ’s are in fact of this form, then a Huffman code does achieve the entropy bound H .

The construction of a Huffman code is best illustrated by example. Imagine a language, Vowellish, with the $N_{ch} = 5$ character alphabet A, E, I, O, and U, occurring with the respective probabilities 0.12, 0.42, 0.09, 0.30, and 0.07. Then the construction of a Huffman code for Vowellish is accomplished in the following table:

Sample page from NUMERICAL RECIPES IN FORTRAN 77: THE ART OF SCIENTIFIC COMPUTING (ISBN 0-521-43064-X)
 Copyright (C) 1986-1992 by Cambridge University Press. Programs Copyright (C) 1986-1992 by Numerical Recipes Software.
 Permission is granted for internet users to make one paper copy for their own personal use. Further reproduction, or any copying of machine-readable files (including this one), to any server computer, is strictly prohibited. To order Numerical Recipes books or CDROMs, visit website <http://www.nr.com> or call 1-800-872-7423 (North America only), or send email to directcustserv@cambridge.org (outside North America).

Node	Stage:	1	2	3	4	5
1	A:	0.12	0.12 ■			
2	E:	0.42	0.42	0.42	0.42 ■	
3	I:	0.09 ■				
4	O:	0.30	0.30	0.30 ■		
5	U:	0.07 ■				
6	UI:		0.16 ■			
7	AUI:			0.28 ■		
8	AUIO:				0.58 ■	
9	EAUIO:					1.00

Here is how it works, proceeding in sequence through N_{ch} stages, represented by the columns of the table. The first stage starts with N_{ch} nodes, one for each letter of the alphabet, containing their respective relative frequencies. At each stage, the two smallest probabilities are found, summed to make a new node, and then dropped from the list of active nodes. (A “block” denotes the stage where a node is dropped.) All active nodes (including the new composite) are then carried over to the next stage (column). In the table, the names assigned to new nodes (e.g., AUI) are inconsequential. In the example shown, it happens that (after stage 1) the two smallest nodes are always an original node and a composite one; this need not be true in general: The two smallest probabilities might be both original nodes, or both composites, or one of each. At the last stage, all nodes will have been collected into one grand composite of total probability 1.

Now, to see the code, you redraw the data in the above table as a tree (Figure 20.4.1). As shown, each node of the tree corresponds to a node (row) in the table, indicated by the integer to its left and probability value to its right. Terminal nodes, so called, are shown as circles; these are single alphabetic characters. The branches of the tree are labeled 0 and 1. The code for a character is the sequence of zeros and ones that lead to it, from the top down. For example, E is simply 0, while U is 1010.

Any string of zeros and ones can now be decoded into an alphabetic sequence. Consider, for example, the string 101111010. Starting at the top of the tree we descend through 1011 to I, the first character. Since we have reached a terminal node, we reset to the top of the tree, next descending through 11 to O. Finally 1010 gives U. The string thus decodes to IOU.

These ideas are embodied in the following routines. Input to the first routine `hufmak` is an integer vector of the frequency of occurrence of the $n_{ch} \equiv N_{ch}$ alphabetic characters, i.e., a set of integers proportional to the p_i 's. `hufmak`, along with `hufapp`, which it calls, performs the construction of the above table, and also the tree of Figure 20.4.1. The routine utilizes a heap structure (see §8.3) for efficiency; for a detailed description, see Sedgewick [7].

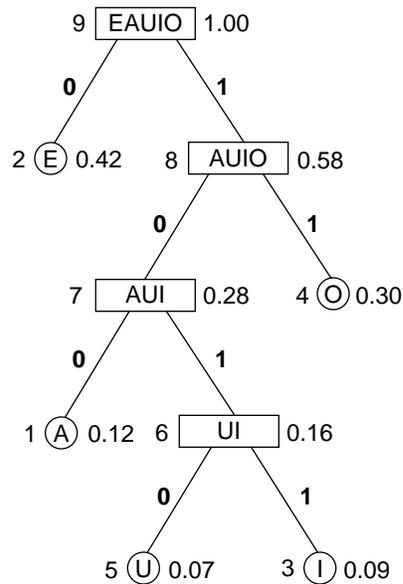


Figure 20.4.1. Huffman code for the fictitious language Vowellish, in tree form. A letter (A, E, I, O, or U) is encoded or decoded by traversing the tree from the top down; the code is the sequence of 0's and 1's on the branches. The value to the right of each node is its probability; to the left, its node number in the accompanying table.

```

SUBROUTINE hufmak(nfreq,nchin,ilong,nlong)
INTEGER ilong,nchin,nlong,nfreq(nchin),MC,MQ
PARAMETER (MC=512,MQ=2*MC-1)
C USES hufapp
    Given the frequency of occurrence table nfreq(1:nchin) of nchin characters, construct
    in the common block /hufcom/ the Huffman code. Returned values ilong and nlong
    are the character number that produced the longest code symbol, and the length of that
    symbol. You should check that nlong is not larger than your machine's word length.
INTEGER ibit,j,k,n,nch,node,nodemx,nused,ibset,index(MQ),
* iup(MQ),icod(MQ),left(MQ),iright(MQ),ncod(MQ),nprob(MQ)
COMMON /hufcom/ icod,ncod,nprob,left,iright,nch,nodemx
SAVE /hufcom/
nch=nchin           Initialization.
nused=0
do 11 j=1,nch
    nprob(j)=nfreq(j)
    icod(j)=0
    ncod(j)=0
    if(nfreq(j).ne.0)then
        nused=nused+1
        index(nused)=j
    endif
enddo 11
do 12 j=nused,1,-1           Sort nprob into a heap structure in index.
    call hufapp(index,nprob,nused,j)
enddo 12
k=nch
1 if(nused.gt.1)then           Combine heap nodes, remarking the heap at each stage.
    node=index(1)
    index(1)=index(nused)
    nused=nused-1
    call hufapp(index,nprob,nused,1)
    k=k+1

```

```

nprob(k)=nprob(index(1))+nprob(node)
left(k)=node           Store left and right children of a node.
iright(k)=index(1)
iup(index(1)) = -k     Indicate whether a node is a left or right child of its parent.
iup(node)=k
index(1)=k
call hufapp(index,nprob,nused,1)
goto 1
endif
nodemx=k
iup(nodemx)=0
do 13 j=1,nch          Make the Huffman code from the tree.
  if(nprob(j).ne.0)then
    n=0
    ibit=0
    node=iup(j)
2    if(node.ne.0)then
      if(node.lt.0)then
        n=ibset(n,ibit)
        node = -node
      endif
      node=iup(node)
      ibit=ibit+1
      goto 2
    endif
    icod(j)=n
    ncod(j)=ibit
  endif
enddo 13
nlong=0
do 14 j=1,nch
  if(ncod(j).gt.nlong)then
    nlong=ncod(j)
    ilong=j-1
  endif
enddo 14
return
END

```

```

SUBROUTINE hufapp(index,nprob,m,1)
INTEGER m,1,MC,MQ
PARAMETER (MC=512,MQ=2*MC-1)
INTEGER index(MQ),nprob(MQ)
  Used by hufmak to maintain a heap structure in the array index(1:1).
INTEGER i,j,k,n
n=m
i=1
k=index(i)
2  if(i.le.n/2)then
    j=i+1
    if (j.lt.n.and.nprob(index(j)).gt.nprob(index(j+1))) j=j+1
    if (nprob(k).le.nprob(index(j))) goto 3
    index(i)=index(j)
    i=j
    goto 2
  endif
3  index(i)=k
return
END

```

Sample page from NUMERICAL RECIPES IN FORTRAN 77: THE ART OF SCIENTIFIC COMPUTING (ISBN 0-521-43064-X)
 Copyright (C) 1986-1992 by Cambridge University Press. Programs Copyright (C) 1986-1992 by Numerical Recipes Software.
 Permission is granted for internet users to make one paper copy for their own personal use. Further reproduction, or any copying of machine-readable files (including this one), to any server computer, is strictly prohibited. To order Numerical Recipes books or CDROMs, visit website <http://www.nr.com> or call 1-800-872-7423 (North America only), or send email to directcustserv@cambridge.org (outside North America).

Once the code is constructed, one encodes a string of characters by repeated calls to `hufenc`, which simply does a table lookup of the code and appends it to the output message.

```

SUBROUTINE hufenc(ich,code,lcode,nb)
INTEGER ich,lcode,nb,MC,MQ
PARAMETER (MC=512,MQ=2*MC-1)
  Huffman encode the single character ich (in the range 0..nch-1), write the result to the
  character array code(1:lcode) starting at bit nb (whose smallest valid value is zero),
  and increment nb appropriately. This routine is called repeatedly to encode consecutive
  characters in a message, but must be preceded by a single initializing call to hufmak.
INTEGER k,l,n,nc,nch,nodemx,ntmp,ibset
INTEGER icod(MQ),left(MQ),iright(MQ),ncod(MQ),nprob(MQ)
LOGICAL btest
CHARACTER*1 code(*)
COMMON /hufcom/ icod,ncod,nprob,left,iright,nch,nodemx
SAVE /hufcom/
k=ich+1          Convert character range 0..nch-1 to array index range 1..nch.
if(k.gt.nch.or.k.lt.1)pause 'ich out of range in hufenc.'
do 11 n=ncod(k),1,-1      Loop over the bits in the stored Huffman code for ich.
  nc=nb/8+1
  if (nc.gt.lcode) pause 'lcode too small in hufenc.'
  l=mod(nb,8)
  if (l.eq.0) code(nc)=char(0)
  if(btest(icod(k),n-1))then Set appropriate bits in code.
    ntmp=ibset(ichar(code(nc)),1)
    code(nc)=char(ntmp)
  endif
  nb=nb+1
enddo 11
return
END

```

Decoding a Huffman-encoded message is slightly more complicated. The coding tree must be traversed from the top down, using up a variable number of bits:

```

SUBROUTINE hufdec(ich,code,lcode,nb)
INTEGER ich,lcode,nb,MC,MQ
PARAMETER (MC=512,MQ=2*MC-1)
  Starting at bit number nb in the character array code(1:lcode), use the Huffman code
  stored in common block /hufcom/ to decode a single character (returned as ich in the
  range 0..nch-1) and increment nb appropriately. Repeated calls, starting with nb = 0
  will return successive characters in a compressed message. The returned value ich=nch
  indicates end-of-message. This routine must be preceded by a single initializing call to
  hufmak.
  Parameters: MC is the maximum value of nch, the input alphabet size.
INTEGER l,nc,nch,node,nodemx
INTEGER icod(MQ),left(MQ),iright(MQ),ncod(MQ),nprob(MQ)
LOGICAL btest
CHARACTER*1 code(lcode)
COMMON /hufcom/ icod,ncod,nprob,left,iright,nch,nodemx
SAVE /hufcom/
node=nodemx          Set node to the top of the decoding tree.
1 continue          Loop until a valid character is obtained.
  nc=nb/8+1
  if (nc.gt.lcode)then Ran out of input; with ich=nch indicating end of message.
    ich=nch
    return
  endif
  l=mod(nb,8)        Now decoding this bit.
  nb=nb+1

```

```

if(bttest(ichar(code(nc)),1))then  Branch left or right in tree, depending on its
    node=iright(node)              value.
else
    node=left(node)
endif
if(node.le.nch)then              If we reach a terminal node, we have a complete character
    ich=node-1                    and can return.
    return
endif
goto 1
END

```

For simplicity, `hufdec` quits when it runs out of code bytes; if your coded message is not an integral number of bytes, and if N_{ch} is less than 256, `hufdec` can return a spurious final character or two, decoded from the spurious trailing bits in your last code byte. If you have independent knowledge of the number of characters sent, you can readily discard these. Otherwise, you can fix this behavior by providing a bit, not byte, count, and modifying the routine accordingly. (When N_{ch} is 256 or larger, `hufdec` will normally run out of code in the middle of a spurious character, and it will be discarded.)

Run-Length Encoding

For the compression of highly correlated bit-streams (for example the black or white values along a facsimile scan line), Huffman compression is often combined with *run-length encoding*: Instead of sending each bit, the input stream is converted to a series of integers indicating how many consecutive bits have the same value. These integers are then Huffman-compressed. The Group 3 CCITT facsimile standard functions in this manner, with a fixed, immutable, Huffman code, optimized for a set of eight standard documents [8,9].

CITED REFERENCES AND FURTHER READING:

- Gallager, R.G. 1968, *Information Theory and Reliable Communication* (New York: Wiley).
- Hamming, R.W. 1980, *Coding and Information Theory* (Englewood Cliffs, NJ: Prentice-Hall).
- Storer, J.A. 1988, *Data Compression: Methods and Theory* (Rockville, MD: Computer Science Press).
- Nelson, M. 1991, *The Data Compression Book* (Redwood City, CA: M&T Books).
- Huffman, D.A. 1952, *Proceedings of the Institute of Radio Engineers*, vol. 40, pp. 1098–1101. [1]
- Ziv, J., and Lempel, A. 1978, *IEEE Transactions on Information Theory*, vol. IT-24, pp. 530–536. [2]
- Cleary, J.G., and Witten, I.H. 1984, *IEEE Transactions on Communications*, vol. COM-32, pp. 396–402. [3]
- Welch, T.A. 1984, *Computer*, vol. 17, no. 6, pp. 8–19. [4]
- Bentley, J.L., Sleator, D.D., Tarjan, R.E., and Wei, V.K. 1986, *Communications of the ACM*, vol. 29, pp. 320–330. [5]
- Jones, D.W. 1988, *Communications of the ACM*, vol. 31, pp. 996–1007. [6]
- Sedgewick, R. 1988, *Algorithms*, 2nd ed. (Reading, MA: Addison-Wesley), Chapter 22. [7]
- Hunter, R., and Robinson, A.H. 1980, *Proceedings of the IEEE*, vol. 68, pp. 854–867. [8]
- Marking, M.P. 1990, *The C Users' Journal*, vol. 8, no. 6, pp. 45–54. [9]

20.5 Arithmetic Coding

We saw in the previous section that a perfect (entropy-bounded) coding scheme would use $L_i = -\log_2 p_i$ bits to encode character i (in the range $1 \leq i \leq N_{ch}$), if p_i is its probability of occurrence. Huffman coding gives a way of rounding the L_i 's to close integer values and constructing a code with those lengths. *Arithmetic coding* [1], which we now discuss, actually does manage to encode characters using noninteger numbers of bits! It also provides a convenient way to output the result not as a stream of bits, but as a stream of symbols in any desired radix. This latter property is particularly useful if you want, e.g., to convert data from bytes (radix 256) to printable ASCII characters (radix 94), or to case-independent alphanumeric sequences containing only A-Z and 0-9 (radix 36).

In arithmetic coding, an input message of any length is represented as a real number R in the range $0 \leq R < 1$. The longer the message, the more precision required of R . This is best illustrated by an example, so let us return to the fictitious language, Vowellish, of the previous section. Recall that Vowellish has a 5 character alphabet (A, E, I, O, U), with occurrence probabilities 0.12, 0.42, 0.09, 0.30, and 0.07, respectively. Figure 20.5.1 shows how a message beginning “IOU” is encoded: The interval $[0, 1)$ is divided into segments corresponding to the 5 alphabetical characters; the length of a segment is the corresponding p_i . We see that the first message character, “I”, narrows the range of R to $0.37 \leq R < 0.46$. This interval is now subdivided into five subintervals, again with lengths proportional to the p_i 's. The second message character, “O”, narrows the range of R to $0.3763 \leq R < 0.4033$. The “U” character further narrows the range to $0.37630 \leq R < 0.37819$. Any value of R in this range can be sent as encoding “IOU”. In particular, the binary fraction .011000001 is in this range, so “IOU” can be sent in 9 bits. (Huffman coding took 10 bits for this example, see §20.4.)

Of course there is the problem of knowing when to stop decoding. The fraction .011000001 represents not simply “IOU,” but “IOU...,” where the ellipses represent an infinite string of successor characters. To resolve this ambiguity, arithmetic coding generally assumes the existence of a special $N_{ch} + 1$ th character, EOM (end of message), which occurs only once at the end of the input. Since EOM has a low probability of occurrence, it gets allocated only a very tiny piece of the number line.

In the above example, we gave R as a binary fraction. We could just as well have output it in any other radix, e.g., base 94 or base 36, whatever is convenient for the anticipated storage or communication channel.

You might wonder how one deals with the seemingly incredible precision required of R for a long message. The answer is that R is never actually represented all at once. At any give stage we have upper and lower bounds for R represented as a finite number of digits in the output radix. As digits of the upper and lower bounds become identical, we can left-shift them away and bring in new digits at the low-significance end. The routines below have a parameter NWK for the number of working digits to keep around. This must be large enough to make the chance of an accidental degeneracy vanishingly small. (The routines signal if a degeneracy ever occurs.) Since the process of discarding old digits and bringing in new ones is performed identically on encoding and decoding, everything stays synchronized.